

By Seokhee Yoon

## COMMENTS ON DATA

Of the many issues related to using statistical data, generalizability is probably the most common question that needs to be addressed for journalists. Called external validity, it is the matter of whether the results of a research project can be extrapolated to the general public. For example, if three-quarters of the respondents in a survey of 100 people believe illegal immigrants are the biggest cause of crime, then can a survey result user say with certainty that three-quarters of the general population believe so, too?

In order to answer this question, two basic areas need to be addressed. One, how was the sample collected? Two, what does the sample consist of?

First of all, because of limited time, money and labor, we cannot afford to conduct a census of every topic we are interested in. Therefore, we use samples to represent the population of interest. However, the quality of the results depends on how the sample was collected. There are two basic kinds of sampling methods: random and non-random. A random sample is a group of subjects (whether they are people, animals, objects, etc) that were extracted from the population in a way that ensured an equal opportunity to be drawn and everyone is drawn by chance. For instance, when you throw a die, each number on the face of the die has an equal chance of showing up and whatever number shows up is due to chance, not manipulation. However, as random samples are not always easy to obtain, researchers sometimes resort to non-random samples such as convenience samples and snowball samples. Using statistics from non-random samples must be approached with caution because it is difficult to generalize the results to the public. For example, if a sample consists of only people who are arrested, then there may be some significant differences between the sample and the general public that affect the results. While 50% of the arrestees may reply that they have used drugs in the past 24 hours, the percentage could be much lower for the general public due to varying life styles and perception about drugs.

In sampling for national telephone surveys, Random Digit Dialing is often employed and the John Jay College survey also used this in their sampling. Random Dialing uses randomly generated telephone numbers for locating samples and although this method is relatively cost effective, simple and convenient to use, it excludes people who do not have telephones. In addition, since many people work outside of homes, they are not there to answer the calls and even when they do answer, they may refuse to participate, especially when the survey is long.

The next question is about the characteristics of the sample. First of all, the size of the sample should be noted. As a general rule, bigger samples are always better and a sample of 1,000 is said to be more generalizable than a sample of 100. Next, the sample should resemble the population. This depends on what the population of interest is. If the population of interest is urban teenagers, then the sample should consist of urban teenagers. If it is voters, like the John Jay survey, then everyone who is selected to participate in it should be verified as registered voters. Also, survey users should be careful not to apply surveys conducted with a sample of limited variability to the general public.

A final tip for reading and incorporating statistics is that even though the sample may have started with 1,000 people, that does not mean all 1,000 people answered all the questions in the survey. The results may vary depending on whether you are looking at valid percentages or total percentages. Valid percentages only take into account valid answers. In other words, they do not count people who refused to answer or people who did not answer properly. If there is a huge chunk of the sample missing in the valid percentages, then it is questionable how externally valid the results are. The results may reflect only certain types of people's answers and therefore, it cannot be generalized to the overall population. Caution is advised before using the statistics as if they are a reflection of the whole sample.